



# A simulation study of maximum likelihood estimation in logistic regression with cured individuals

Aba Diop, Aliou Diop, Jean-François Dupuy

## ► To cite this version:

Aba Diop, Aliou Diop, Jean-François Dupuy. A simulation study of maximum likelihood estimation in logistic regression with cured individuals. 2011. hal-00636486

**HAL Id: hal-00636486**

**<https://hal.science/hal-00636486>**

Submitted on 27 Oct 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A simulation study of maximum likelihood estimation in logistic regression with cured individuals

Aba DIOP

*MIA, Université de La Rochelle, France and LERSTAD, Université Gaston Berger, Saint Louis, Sénégal.*

*Email: aba.diop@univ-lr.fr*

Aliou DIOP

*LERSTAD, Université Gaston Berger, Saint Louis, Sénégal.*

*Email: aliou.diop@ugb.edu.sn*

Jean-François DUPUY†

*IRMAR-Institut National des Sciences Appliquées de Rennes, France.*

*Email: Jean-Francois.Dupuy@insa-rennes.fr*

**Abstract.** The logistic regression model is widely used to investigate the relationship between a binary outcome  $Y$  and a set of potential predictors  $\mathbf{X}$ . Diop *et al.* (2011) present some conditions under which the maximum likelihood estimator is consistent and asymptotically normal in the logistic regression model with a cure fraction. So far, however, only limited simulation results are available to judge the quality of this estimator in finite samples. Therefore in this paper, we conduct a detailed simulation study of its numerical properties. We evaluate its accuracy and the quality of the normal approximation of its asymptotic distribution. We also study the quality of the approximation for constructing asymptotic Wald-type tests of hypothesis. Finally, we consider the problem of estimating the conditional probability of the outcome. Our results indicate that when the proportion of cured individuals is moderate to moderately large, and the sample size is large enough, reliable statistical inferences can be obtained for the regression effects and the probability of the outcome. Our results also indicate that the approximations can be problematic when the cure fraction is very large.

**Keywords:** Zero-inflation, mixture model, maximum likelihood estimation

## 1. Introduction

Logistic regression has become a standard tool to investigate the relationship between a binary response  $Y$  and a set of potential predictors  $\mathbf{X}$ . In the medical setting for example, the response may represent the infection status with respect to some disease. If  $Y_i$  denotes the infection status for the  $i$ -th individual in a sample of size  $n$  ( $Y_i = 1$  if the individual is infected, and  $Y_i = 0$  otherwise) and  $\mathbf{X}_i = (1, X_{i2}, \dots, X_{ip})'$  is the corresponding predictor, logistic regression models the conditional probability  $\mathbb{P}(Y_i = 1|\mathbf{X}_i)$  of infection as

$$\log \left( \frac{\mathbb{P}(Y_i = 1|\mathbf{X}_i)}{1 - \mathbb{P}(Y_i = 1|\mathbf{X}_i)} \right) = \beta' \mathbf{X}_i, \quad (1)$$

†Corresponding author

where  $\beta = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$  is an unknown regression parameter to be estimated. Estimation and testing procedures in the model (1) are well established (see, for example, Hosmer and Lemeshow (2000) and Hilbe (2009)) and are available in standard statistical softwares. In particular, the maximum likelihood estimator (MLE) of  $\beta$  was shown to be consistent, and approximately normally distributed in large samples (Gouriéroux and Monfort (1981), Fahrmeir and Kaufmann (1985)).

Recently, Diop *et al.* (2011) have considered the problem of estimation in the model (1) when there is a cured fraction in the sample. In medical studies, it often arises that a proportion of the study subjects cannot experience the outcome of interest. Such individuals are said to be cured, or immune. The population under study can then be considered as a mixture of cured and susceptible subjects, where a subject is said to be susceptible if he would eventually experience the outcome of interest. One problem arising in this setting is that it is usually unknown who are the susceptible, and the cured subjects (unless the outcome of interest has been observed). Consider, for example, the occurrence of infection from some disease to be the outcome of interest. Then, if a subject is uninfected, the investigator usually does not know whether this subject is immune to the infection, or susceptible albeit still uninfected.

Estimating a regression model with a cure fraction can be viewed as a zero-inflated regression problem. Zero-inflation occurs in the analysis of count data when the observations contain more zeros than expected. Failure to account for these extra zeros is known to result in biased parameter estimates and inferences. Motivated by various applications in public health, epidemiology, sociology, engineering, agriculture, a variety of zero-inflated regression models have been developed and extended, such as the zero-inflated Poisson (ZIP) model (see, among others, Lambert (1992), Dietz and Böhning (2000), Ridout *et al.* (2001), Lam *et al.* (2006), Xiang *et al.* (2007)), the zero-inflated binomial (ZIB) model (see Hall (2000)), the zero-inflated negative binomial (ZINB) model (see Moghimbeigi *et al.* (2008)). Various other models and numerous references can be found in Famoye and Singh (2006), Lee *et al.* (2006), Kelley and Anderson (2008), and Moghimbeigi *et al.* (2009).

By assuming a logistic regression model for the probability of being cured, Diop *et al.* (2011) developed a new zero-inflated Bernoulli regression model with logit links for both the binary response of interest (the probability of infection, say), and the zero-inflation probability (of being cured). The authors investigated, mainly theoretically, the issues of identifiability and estimation in this model. In particular, they proposed a set of sufficient conditions for model identifiability, and they proved the consistency and asymptotic normality of the MLE.

In this paper, as a supplement of the theoretical work of Diop *et al.* (2011), we conduct a detailed simulation study of the numerical behavior of their estimator. Precisely, we evaluate the influence of various factors (sample size, proportion of immunes in the sample, proportion of infected among the susceptibles) on the accuracy of the estimator and on the quality of the Gaussian approximation of its asymptotic distribution. The performance (in terms of power and level) of the Wald-type test of " $\beta = 0$ " is also investigated. We also consider the problem of estimating the probability of infection  $p(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$ , for some value  $\mathbf{x}$  of the covariates.

The rest of the paper is organized as follows. In the Section 2, we describe the model and the estimation procedure proposed by Diop *et al.* (2011), and we recall the theoretical properties of the resulting estimator of the regression parameter  $\beta$ . We also construct asymptotic confidence intervals for a probability of infection  $p(\mathbf{x})$ . Section 3 describes the detailed simulation study. A discussion and some conclusions and perspectives are given in

the Section 4.

## 2. Model and estimation

Let  $\mathcal{O}_i = (Y_i, S_i, \mathbf{X}_i, \mathbf{Z}_i)$ ,  $i = 1, \dots, n$  be independent copies of the random vector  $\mathcal{O} = (Y, S, \mathbf{X}, \mathbf{Z})$ , where for every  $i$ ,  $Y_i$  is a binary response indicating the infection status (with respect to some disease) of the  $i$ -th individual (that is,  $Y_i = 1$  if the individual  $i$  is infected, and  $Y_i = 0$  otherwise), and  $S_i$  is a binary variable indicating whether the individual is susceptible to the infection ( $S_i = 1$ ) or immune ( $S_i = 0$ ). If  $Y_i = 0$ , then the value of  $S_i$  is unknown.  $\mathbf{X}_i = (1, X_{i2}, \dots, X_{ip})'$  and  $\mathbf{Z}_i = (1, Z_{i2}, \dots, Z_{iq})'$  are covariate vectors related to the infection status and immunity respectively.  $\mathbf{X}_i$  and  $\mathbf{Z}_i$  may contain quantitative and qualitative components, and may even share some components.

The zero-inflated Bernoulli regression model described by Diop *et al.* (2011) is defined by the following equations for the infection status:

$$\begin{cases} \log \left( \frac{\mathbb{P}(Y=1|\mathbf{X}_i, S_i)}{1 - \mathbb{P}(Y=1|\mathbf{X}_i, S_i)} \right) = \beta' \mathbf{X}_i & \text{if } \{S_i = 1\} \\ \mathbb{P}(Y = 1|\mathbf{X}_i, S_i) = 0 & \text{if } \{S_i = 0\} \end{cases} \quad (2)$$

and by the following model for immunity:

$$\log \left( \frac{\mathbb{P}(S = 1|\mathbf{Z}_i)}{1 - \mathbb{P}(S = 1|\mathbf{Z}_i)} \right) = \theta' \mathbf{Z}_i. \quad (3)$$

In this model,  $\beta = (\beta_1, \dots, \beta_p)' \in \mathbb{R}^p$  is an unknown regression parameter of interest ( $\beta$  measures the association between the potential predictors  $\mathbf{X}_i$  and the risk of infection for a susceptible individual), and  $\theta = (\theta_1, \dots, \theta_q)' \in \mathbb{R}^q$  is an unknown nuisance parameter. Letting  $\psi := (\beta', \theta')'$ , the log-likelihood for  $\psi$  from the sample  $\mathcal{O}_1, \dots, \mathcal{O}_n$  (where  $S_i$  is unknown when  $Y_i = 0$ ) is

$$\begin{aligned} l_n(\psi) &= \sum_{i=1}^n \left\{ Y_i(\beta' \mathbf{X}_i + \theta' \mathbf{Z}_i) + (1 - Y_i) \log(1 + e^{\beta' \mathbf{X}_i} + e^{\theta' \mathbf{Z}_i}) \right. \\ &\quad \left. - \log(1 + e^{\beta' \mathbf{X}_i}) - \log(1 + e^{\theta' \mathbf{Z}_i}) \right\} \\ &:= \sum_{i=1}^n l(\psi; \mathcal{O}_i). \end{aligned}$$

The MLE  $\hat{\psi}_n := (\hat{\beta}'_n, \hat{\theta}'_n)'$  of  $\psi$  is defined as the solution of the score equation  $\partial l_n(\psi) / \partial \psi = 0$ , which can be solved, for example, using the `optim` function of the software R.

The main assumptions underlying the asymptotic results proved by Diop *et al.* (2011) are the following:

- A1** The covariates are bounded. For every  $i = 1, 2, \dots$ ,  $j = 2, \dots, p$ ,  $k = 2, \dots, q$ ,  $\text{var}[X_{ij}] > 0$  and  $\text{var}[Z_{ik}] > 0$ . For every  $i = 1, 2, \dots$ , the  $X_{ij}$  ( $j = 1, \dots, p$ ) are linearly independent, and the  $Z_{ik}$  ( $k = 1, \dots, q$ ) are linearly independent.
- A2** There exists a continuous covariate  $V$  which is in  $\mathbf{X}$  but not in  $\mathbf{Z}$  that is, if  $\beta_V$  and  $\theta_V$  denote the coefficients of  $V$  in the linear predictors (2) and (3) respectively, then  $\beta_V \neq 0$  and  $\theta_V = 0$ . At a model-building stage, it is known that  $V$  is in  $\mathbf{X}$ .

**Remark 1.** The condition A1 is a classical condition for identifiability and asymptotic results in standard logistic regression (see, for example, Gouriéroux and Monfort (1981)). The condition A2, which imposes some restrictions on the covariates, is required for the identifiability of  $\psi$  in the joint model (2)-(3) (we may alternatively assume that the continuous covariate  $V$  is in  $\mathbf{Z}$  but not in  $\mathbf{X}$ ). In the following, we will assume that  $V$  is in  $\mathbf{X}$  but not in  $\mathbf{Z}$ , with  $\beta_V := \beta_l$  for some  $l \in \{2, \dots, p\}$ , and for the  $i$ -th individual, we will denote  $V_i$  by  $X_{il}$ . We refer to Diop *et al.* (2011) for a detailed discussion about this condition.

Let  $\mathcal{I}_\psi = -\mathbb{E}[\partial^2 l(\psi; \mathcal{O}) / \partial \psi \partial \psi']$ . Under the conditions stated above, Diop *et al.* (2011) prove the following result:

**Theorem.** *As  $n$  tends to infinity,  $\sqrt{n}(\hat{\psi}_n - \psi)$  converges in distribution to a zero-mean Gaussian vector, whose covariance matrix  $\mathcal{I}_\psi^{-1}$  can be consistently estimated by  $\hat{\mathcal{I}}_{\hat{\psi}_n}^{-1}$ , where  $\hat{\mathcal{I}}_{\hat{\psi}_n} = -n^{-1}(\partial^2 l_n(\psi) / \partial \psi \partial \psi')|_{\psi=\hat{\psi}_n}$ .*

The parameter of interest in the model (2)-(3) is  $\beta$ . The asymptotic distribution of  $\hat{\beta}_n$  now easily follows:

**Corollary 1.** *Let  $M$  be the  $(p \times (p+q))$  block-matrix  $[I_p, 0_{p,q}]$ , where  $I_p$  denotes the identity matrix of order  $p$  and  $0_{p,q}$  is the  $(p \times q)$  matrix whose components are all equal to 0. Then  $\sqrt{n}(\hat{\beta}_n - \beta)$  converges in distribution to a zero-mean Gaussian vector with covariance matrix  $M\mathcal{I}_\psi^{-1}M'$ , which is the upper-left  $(p \times p)$  block of  $\mathcal{I}_\psi^{-1}$ .*

The convergence in distribution of  $\hat{\beta}_n$  can now be used to make statistical inference about  $\beta$ . For example, if one wishes to test the null hypothesis  $H_0 : \beta_l = 0$  against the alternative  $H_1 : \beta_l \neq 0$  (for some  $1 \leq l \leq p$ ), one can rely on the usual Wald-type test, which rejects  $H_0$  at the asymptotic level  $\alpha$  ( $0 < \alpha < 1$ ) if

$$\left| \frac{\hat{\beta}_{n,l}}{\sqrt{\frac{(M\hat{\mathcal{I}}_{\hat{\psi}_n}^{-1}M')_{ll}}{n}}} \right| > u_{1-\frac{\alpha}{2}},$$

where  $u_{1-\frac{\alpha}{2}}$  is the quantile of order  $(1 - \frac{\alpha}{2})$  of the standard normal law,  $\hat{\beta}_{n,l}$  is the  $l$ -th component of  $\hat{\beta}_n$ , and  $(M\hat{\mathcal{I}}_{\hat{\psi}_n}^{-1}M')_{ll}$  denotes the  $l$ -th diagonal component of  $M\hat{\mathcal{I}}_{\hat{\psi}_n}^{-1}M'$ . This corollary can also be used to estimate a given probability  $p(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$  of infection (this result was not mentioned by Diop *et al.* (2011)):

**Corollary 2.** *Let  $\mathbf{x}$  be a given value of the covariate  $\mathbf{X}$ . As  $n$  tends to infinity,  $\sqrt{n}(\hat{\beta}_n - \beta)' \mathbf{x}$  converges in distribution to a zero-mean Gaussian random variable with variance  $\mathbf{x}' M \mathcal{I}_\psi^{-1} M' \mathbf{x}$ . An asymptotic  $(1 - \alpha)$ -level confidence interval for  $p(\mathbf{x})$  is therefore*

$$\left[ \frac{e^{u_n(\mathbf{x})}}{1 + e^{u_n(\mathbf{x})}}, \frac{e^{v_n(\mathbf{x})}}{1 + e^{v_n(\mathbf{x})}} \right],$$

where

$$u_n(\mathbf{x}) = \hat{\beta}'_n \mathbf{x} - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\mathbf{x}' M \hat{\mathcal{I}}_{\hat{\psi}_n}^{-1} M' \mathbf{x}}{n}} \quad \text{and} \quad v_n(\mathbf{x}) = \hat{\beta}'_n \mathbf{x} + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\mathbf{x}' M \hat{\mathcal{I}}_{\hat{\psi}_n}^{-1} M' \mathbf{x}}{n}}.$$

In a limited simulation study, Diop *et al.* (2011) have investigated the numerical behavior of  $\hat{\beta}_n$  for some very simple models. In the present paper, we undertake a much more detailed numerical analysis of this estimator, focusing on various measures of its accuracy (namely the bias, and the root mean-square and mean absolute errors), on the quality of the Gaussian approximation of its asymptotic distribution, and on the accuracy of an estimator of the probability of infection  $p(\mathbf{x})$ , for some  $\mathbf{x}$ .

### 3. Simulation study

#### 3.1. Study design

The simulation setting is as follows. We consider the following models for the infection status:

$$\begin{cases} \log \left( \frac{\mathbb{P}(Y=1|\mathbf{X}_i, S_i)}{1 - \mathbb{P}(Y=1|\mathbf{X}_i, S_i)} \right) = \beta_1 + \beta_2 X_{i2} + \beta_3 Z_{i2} + \beta_4 Z_{i3} + \beta_5 Z_{i4} & \text{if } S_i = 1 \\ \mathbb{P}(Y = 1|\mathbf{X}_i, S_i) = 0 & \text{if } S_i = 0 \end{cases} \quad (4)$$

and the immunity status:

$$\log \left( \frac{\mathbb{P}(S = 1|\mathbf{Z}_i)}{1 - \mathbb{P}(S = 1|\mathbf{Z}_i)} \right) = \theta_1 + \theta_2 Z_{i2} + \theta_3 Z_{i3} + \theta_4 Z_{i4}, \quad (5)$$

where  $X_{i2} \sim N(0, 1)$ ,  $Z_{i2} \sim N(1, 1)$ , and  $Z_{i3}$  and  $Z_{i4}$  are indicator variables built from a categorical variable with 3 categories. Note that  $X_2$  plays the role of the continuous covariate  $V$  in the condition A2. Note also that the models for infection and immunity share no less than three covariates, including both continuous and discrete variables (this is the least favorable case with respect to the identifiability of the parameters). An i.i.d. sample of size  $n$  of the vector  $(Y, S, \mathbf{X}, \mathbf{Z})$  is generated from this model, and for each individual  $i$ , we get a realization  $(y_i, s_i, \mathbf{x}_i, \mathbf{z}_i)$ , where  $s_i$  is considered as unknown if  $y_i = 0$ . The maximum likelihood estimator  $\hat{\beta}_n$  of  $\beta = (\beta_1, \beta_2, \beta_3, \beta_4, \beta_5)'$  is obtained from this incomplete dataset by solving the score equation given in Section 2, using the `optim` function in **R**. As a by-product of the method, an estimate is also obtained for the nuisance parameter  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)'$ .

The finite-sample behavior of the maximum likelihood estimator  $\hat{\beta}_n$  is assessed for several sample sizes ( $n = 100, 500, 1000, 1500$ ) and various values for the percentage of immunes in the sample, namely 25%, 50%, and 75%. We also consider different values for the proportion of infected individuals among the susceptibles (30% and 70%). The desired proportions of immunes and infected are obtained by choosing appropriate values for  $\beta$  and  $\theta$ . The following values are considered for  $\beta$ :

- model  $\mathcal{M}_1$ :  $\beta = (-1.7, -2, -3.4, 5, 0.3)'$ . Using these values, approximately 30% of the susceptibles are infected.
- model  $\mathcal{M}_2$ :  $\beta = (1.5, -2.3, 2.5, -3.5, 0.5)'$ . Approximately 70% of the susceptibles are infected.
- model  $\mathcal{M}_3$ :  $\beta = (-1.7, -2.8, 0, -0.7, 1.1)'$ . Approximately 30% of the susceptibles are infected.

- model  $\mathcal{M}_4$ :  $\beta = (1.5, -2, 0, 3.5, -4)'$ . Approximately 70% of the susceptibles are infected.

Letting  $\beta_3 = 0$  (models  $\mathcal{M}_3$  and  $\mathcal{M}_4$ ) implies that  $Z_2$  does not influence the risk of infection. However, when estimating the model (4)-(5) using Diop *et al.* (2011)'s procedure, we use the whole set of covariates  $(X_2, Z_2, Z_3, Z_4)$ . This allows us to evaluate the level of a Wald-type test of nullity of  $\beta_3$ .

### 3.2. Results

For each configuration **sample size**  $\times$  **percentage of immunes**  $\times$  **percentage of infected among susceptibles** of the design parameters,  $N = 1500$  samples are obtained. Based on these  $N$  repetitions, we obtain averaged values for the estimates of the  $\beta_l$  ( $l = 1, \dots, 5$ ), which are calculated as  $N^{-1} \sum_{j=1}^N \hat{\beta}_{l,n}^{(j)}$ , where  $\hat{\beta}_n^{(j)} = (\hat{\beta}_{1,n}^{(j)}, \dots, \hat{\beta}_{5,n}^{(j)})'$  is the estimate obtained from the  $j$ -th simulated sample. For each of the parameters  $\beta_l$ , we also obtain the empirical root mean square (RMSE) and mean absolute errors (MAE), based on the  $N$  samples. Similar results are obtained for the nuisance parameter  $\theta$ . When  $\beta_l \neq 0$  (respectively  $\beta_l = 0$ ), we obtain the empirical power (respectively the empirical size) of the Wald test at the 5% level for testing  $H_0 : \beta_l = 0$ . The null hypothesis  $H_0 : \beta_l = 0$  is the hypothesis that the predictor  $X_l$  does not influence the risk of infection of susceptible individuals. Tables 1 and 2 give the results for the models  $\mathcal{M}_1$  and  $\mathcal{M}_2$  respectively. There,  $(\cdot)$  indicates the root mean square error,  $[\cdot]$  indicates the mean absolute error, and  $*$  is the empirical power of the Wald test.

**Table 1 about here**

**Table 2 about here**

In order to investigate the level of Wald-type tests of hypothesis of the form  $H_0 : \beta_l = 0$ , we simulate samples from the models  $\mathcal{M}_3$  and  $\mathcal{M}_4$ . The results are given in the tables 3 and 4. In these tables,  $*$  indicates the empirical power of the Wald test (when  $\beta_l \neq 0$ ) and  $\dagger$  indicates the empirical size (when  $\beta_l = 0$ ).

**Table 3 about here**

**Table 4 about here**

If there were no immunes in the sample, a usual logistic regression model could be fitted to the data. The results for such an hypothetical case are interesting since they provide a benchmark for evaluating the performance of the proposed estimator. Tables 5 and 6 provide the results for the models  $\mathcal{M}_1$  and  $\mathcal{M}_3$  respectively, if there were no immunes in the samples. The results for the other models yield similar conclusions and are therefore omitted.

**Table 5 about here**

**Table 6 about here**

From all these results, it appears that the maximum likelihood estimator proposed by Diop *et al.* (2011) provides a reasonable approximation of the true regression parameter, even

when the immune fraction is high. Overall, the bias of  $\hat{\beta}_n$  stays limited while its variability increases with the percentage of immunes. This increase is particularly noticeable when the sample size is small ( $n = 100$ ) and when the sample size is moderate ( $n = 500$ ) with a high percentage of immunes (75%). As a consequence, the power of Wald tests of nullity of regression coefficients can be low (compared to the case where there are no immunes) when the percentage of immunes is very high and/or the sample size is small. For moderately large to large sample sizes ( $n \geq 500$ , say), the power of the Wald test indicates good numerical performance of the inferential procedure proposed by Diop *et al.* (2011), provided that the magnitude of the regression coefficient is sufficiently large. The level of the Wald test is globally respected when the percentage of immunes is moderate (25%) and degrades when this percentage increases.

Then, we compare these results to the ones obtained from a "naive" method where: i) we consider every individual  $i$  such that  $\{Y_i = 0\}$  as being susceptible but uninfected (that is, we ignore the eventual immunity of this individual) and ii) we apply a usual logistic regression analysis to the resulting dataset. The results of such "naive" analysis for the model  $\mathcal{M}_1$  are given in the table 7. The results for the other models yield similar observations and thus, they are not given here.

#### Table 7 about here

From this table, it appears that missing the immunity present in the sample results in strongly biased estimates of the true association between the covariates and the risk of infection. The variability of the estimates is also very important, even for very large sample sizes.

In logistic regression, it is usually of interest to estimate the probability of infection  $p(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$ , for some given value  $\mathbf{x}$  of the covariates. An obvious estimate of  $p(\mathbf{x})$  is  $\hat{p}_n(\mathbf{x}) := \exp(\hat{\beta}'_n \mathbf{x}) / (1 + \exp(\hat{\beta}'_n \mathbf{x}))$ . In table 8, we investigate the numerical properties of this estimator (restraining ourselves to one value of  $p(\mathbf{x})$  for each of the models  $\mathcal{M}_1$  and  $\mathcal{M}_2$ ). For every combination of the simulation design parameters, we obtain the average estimated probability and the corresponding RMSE and MAE. We also provide the empirical coverage probability of asymptotic 95%-confidence intervals for  $p(\mathbf{x})$ , and the average length of these intervals.

#### Table 8 about here

From this table, it appears that  $\hat{p}_n(\mathbf{x})$  provides a reasonable estimate of  $p(\mathbf{x})$  even when the immune fraction is high. The coverage probabilities and average interval lengths indicate good performance (compared to the no-immunes case) of the confidence intervals for  $p(\mathbf{x})$  obtained from Diop *et al.* (2011)'s results, when the cure fraction is moderate (even for  $n = 100$ ), and when the cure fraction is moderately large (50%) and the sample size is large ( $n \geq 1000$ ).

We also investigate the quality of the normal approximation of the large-sample distribution of  $\hat{\beta}_n$ . For each configuration of the design parameters, we obtain histograms of the  $\hat{\beta}_{n,l}^{(j)}$  ( $j = 1, \dots, N$ ), along with the corresponding Q-Q plots. The results are provided for the coefficients  $\beta_2$  and  $\beta_3$  of the model  $\mathcal{M}_1$ . The results for the other parameters and models yield similar observations and are therefore omitted.

#### Figures 1 to 8 about here



From these figures, the normal approximation seems to be reasonably satisfied when the immune fraction is moderate (25%) and the sample size is sufficiently large ( $n \geq 500$ ). When the proportion of immunes attains 50%, a larger sample size ( $n \geq 1000$ , say) is required to consider this approximation as being still reasonable. When the immune proportion is very large (75%), the distribution of the estimator can be highly skewed, particularly when the sample size is small.

Finally, we investigate the quality of the normal approximation of the distribution of  $\hat{p}_n(\mathbf{x})$ . Histograms and Q-Q plots of the  $\hat{p}_n^{(j)}(\mathbf{x})$  ( $j = 1, \dots, N$ ) are obtained for one value of  $p(\mathbf{x})$  in the model  $\mathcal{M}_1$ , and are given in Figures 9-12.

#### Figures 9 to 12 about here

From these figures, the normal approximation is reasonably satisfied when the immune fraction is moderate to moderately large (25%-50%) provided that the sample size is sufficiently large ( $n \geq 1000$ , say). The distribution of  $\hat{p}_n(\mathbf{x})$  can be highly skewed otherwise.

Overall, these results indicate that a reliable statistical inference on the regression effects and probabilities of event in the logistic regression model with a cure fraction should be based on a sample having, at least, a moderately large size ( $n \geq 500$ , say) when the immune fraction is low (25%), or a large size ( $n \geq 1000$ ) when the immunity attains 50% of the sample. When the immune proportion is very large (75%), the results should be considered carefully, considering the increase in the variability of the estimates and the skewness of their distributions.

## 4. Discussion

In this paper, we have conducted a detailed numerical investigation of the maximum likelihood estimators proposed by Diop *et al.* (2011) for estimating the logistic regression model when there is a cure fraction in the sample. This problem can also be viewed as a problem of zero-inflation in the logistic regression model, or equivalently, as a problem of statistical inference in logistic regression from a mixture of cured and susceptible individuals. From our investigations, the maximum likelihood estimators of both the regression effects of interest and probabilities of event (such as infection, say) perform quite well under appropriate conditions regarding the sample size and immune proportion. Overall, reliable statistical inferences (point estimation, confidence intervals, hypothesis testing) should be obtained when the immune fraction is moderate (25%) and the sample size is at least 500, or the immune fraction is moderately large (50%) and the sample size is at least 1000.

Several open problems still deserve attention. First, our results were obtained under the assumption that the cure model is correctly specified. It is now of interest to investigate the effect on the whole inference of a misspecification of this model. One may also consider estimating the cure model by using more flexible methods, such as nonparametric kernel regression. Another stimulating topic of interest is to consider the problem of immunes in the more general semiparametric logistic regression model. Both topics are the subject for our future research.

In logistic regression, it is sometimes of interest to make inference about the probability of event  $\mathbb{P}(Y = 1 | \mathbf{X}, S = 1)$  across the whole range of the predictors  $\mathbf{X}$ . Pointwise confidence intervals are not adequate for that purpose. The calculation of simultaneous confidence bands for the probabilities  $\{p(\mathbf{x}), \mathbf{x} \in \mathcal{X}\}$  thus constitutes another issue of (both methodological and practical) interest. Several methods have been proposed for constructing confidence bands in regression models (see, for example, Brand *et al.* (1973) and Hauck

(1983) for the standard logistic regression model). Extending and evaluating these methods in the case of logistic regression with a cure fraction constitutes a stimulating topic for future research.

## Acknowledgements

This research was supported by AIRES-Sud (AIRES-Sud is a program from the French Ministry of Foreign and European Affairs, implemented by the "Institut de Recherche pour le Développement", IRD-DSF), by the "Service de Coopération et d'Action Culturelle" of the French Embassy in Senegal, and by Edulink (program 9-ACP-RPR-118#18). The authors also acknowledge grants from the "Ministère de la Recherche Scientifique" of Senegal.

## References

- Brand, R., Pinnock, D., and Jackson, K. (1973). Large sample confidence bands for the logistic response curve and its inverse. *American Statistician* **27**, 157–160.
- Dietz, E. and Böhning, D., 2000. On estimation of the Poisson parameter in zero-modified Poisson models. *Computational Statistics & Data Analysis* **34**, 441–459.
- Diop, A., Diop, A., and Dupuy, J.-F. (2011). Maximum likelihood estimation in the logistic regression model with a cure fraction. *Electronic Journal of Statistics* **5**, 460–483.
- Fahrmeir, L. and Kaufmann, H. (1985). Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *The Annals of Statistics* **13**, 342–368.
- Famoye, F. and Singh, K.P., 2006. Zero-inflated generalized Poisson regression model with an application to domestic violence data. *Journal of Data Science* **4**, 117–130.
- Gouriéroux, C. and Monfort, A., 1981. Asymptotic properties of the maximum likelihood estimator in dichotomous logit models. *Journal of Econometrics* **17**, 83–97.
- Hall, D.B., 2000. Zero-inflated Poisson and binomial regression with random effects: a case study. *Biometrics* **56**, 1030–1039.
- Hauck, W.W. (1983). A note on confidence bands for the logistic response curve. *American Statistician* **37**, 158–160.
- Hilbe, J.M., 2009. *Logistic regression models*. Chapman & Hall: Boca Raton.
- Hosmer, D.W. and Lemeshow, S., 2000. *Applied logistic regression*. Wiley: New York.
- Kelley, M. E. and Anderson, S.J., 2008. Zero inflation in ordinal data: incorporating susceptibility to response through the use of a mixture model. *Statistics in Medicine* **27**, 3674–3688.
- Lam, K.F., Xue, H. and Cheung, Y.B., 2006. Semiparametric analysis of zero-inflated count data. *Biometrics* **62**, 996–1003.
- Lambert, D., 1992. Zero-inflated Poisson regression, with an application to defects in manufacturing. *Technometrics* **34**, 1–14.

- Lee, A.H., Wang, K., Scott, J.A., Yau, K.K.W. and McLachlan, G. J., 2006. Multi-level zero-inflated Poisson regression modelling of correlated count data with excess zeros. *Statistical Methods in Medical Research* **15**, 47-61.
- Moghimbeigi, A., Eshraghian, M.R., Mohammad, K., and McArdle, B., 2009. Multilevel zero-inflated negative binomial regression modeling for over-dispersed count data with extra zeros. *Journal of Applied Statistics* **35**, 1193–1202.
- Moghimbeigi, A., Eshraghian, M.R., Mohammad, K., and McArdle, B., 2009. A score test for zero-inflation in multilevel count data. *Computational Statistics & Data Analysis* **53**, 1239–1248.
- Ridout, M., Hinde, J., and Demétrio, C.G.B., 2001. A score test for testing a zero-inflated Poisson regression model against zero-inflated negative binomial alternatives. *Biometrics* **57**, 219–223.
- Xiang, L., Lee, A. H., Yau, K. K. W., and McLachlan, G. J., 2007. A score test for overdispersion in zero-inflated Poisson mixed regression model. *Statistics in Medicine* **26**, 1608–1622.

**Table 1. Simulation results for model  $\mathcal{M}_1$ :**  $\beta = (-1.7, -2, -3.4, 5, .3)$  (percentage of infected among the susceptibles: 30%).

n	$\hat{\beta}_n$					$\hat{\theta}_n$			
	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$	$\hat{\theta}_{1,n}$	$\hat{\theta}_{2,n}$	$\hat{\theta}_{3,n}$	$\hat{\theta}_{4,n}$
<b>Percentage of immunes = 25%, <math>\theta = (.71, 1, 2, -3)</math></b>									
100	-1.688 (1.625) [1.330]	-2.019 (1.160) [0.689] 0.631*	-3.417 (1.704) [1.387] 0.539*	4.648 (1.866) [1.525] 0.316*	0.322 (1.880) [1.524] 0.013*	0.667 (2.660) [2.199]	1.154 (1.462) [1.161]	2.153 (1.715) [1.410]	-3.088 (2.663) [2.156]
500	-1.710 (0.966) [0.734]	-2.006 (0.427) [0.330] 0.997*	-3.392 (0.923) [0.686] 0.993*	4.991 (1.167) [0.907] 0.984*	0.311 (1.581) [1.329] 0.090*	0.677 (0.849) [0.711]	1.076 (0.804) [0.643]	2.079 (1.144) [0.910]	-2.994 (2.090) [1.663]
1000	-1.702 (0.579) [0.456]	-2.004 (0.297) [0.233] 1*	-3.398 (0.584) [0.412] 0.998*	4.968 (0.797) [0.612] 1*	0.305 (0.843) [0.716] 0.107*	0.697 (0.749) [0.611]	1.046 (0.623) [0.477]	2.026 (0.779) [0.605]	-2.997 (1.127) [0.910]
1500	-1.720 (0.492) [0.384]	-1.998 (0.272) [0.206] 1*	-3.410 (0.474) [0.332] 1*	4.971 (0.649) [0.484] 1*	0.305 (0.794) [0.675] 0.095*	0.709 (0.607) [0.487]	1.035 (0.475) [0.361]	2.013 (0.614) [0.484]	-3.002 (0.979) [0.778]
<b>Percentage of immunes = 50%, <math>\theta = (-.3, -1, 2.1, 1)</math></b>									
100	-1.767 (2.068) [1.682]	-2.105 (1.013) [0.784] 0.335*	-3.341 (1.783) [1.464] 0.337*	5.334 (2.557) [2.150] 0.127*	0.317 (2.758) [2.312] 0*	-0.377 (2.672) [2.141]	-1.123 (1.965) [1.265]	2.212 (2.156) [1.760]	1.090 (2.851) [2.366]
500	-1.716 (1.417) [1.085]	-2.081 (0.545) [0.452] 1*	-3.576 (0.930) [0.743] 1*	5.217 (1.704) [1.383] 0.839*	0.294 (1.704) [1.320] 0.057*	-0.342 (1.079) [0.875]	-1.092 (0.760) [0.554]	2.097 (1.518) [1.177]	1.078 (1.746) [1.401]
1000	-1.701 (0.780) [0.650]	-2.076 (0.391) [0.316] 1*	-3.529 (0.627) [0.495] 1*	5.015 (1.072) [0.866] 0.984*	0.304 (1.113) [0.892] 0.045*	-0.318 (0.743) [0.607]	-1.056 (0.473) [0.352]	2.105 (1.028) [0.770]	1.031 (1.160) [0.905]
1500	-1.698 (0.694) [0.568]	-2.013 (0.294) [0.242] 1*	-3.482 (0.489) [0.381] 1*	5.007 (0.857) [0.694] 0.999*	0.303 (0.865) [0.685] 0.057*	-0.315 (0.609) [0.497]	-1.023 (0.384) [0.290]	2.103 (0.885) [0.642]	1.024 (0.926) [0.721]
<b>Percentage of immunes = 75%, <math>\theta = (.4, -1, -.6, -2)</math></b>									
100	-1.661 (2.139) [1.754]	-2.131 (2.127) [1.720] 0.043*	-3.387 (2.803) [2.394] 0.064*	4.830 (3.283) [2.849] 0.005*	0.332 (3.661) [2.811] 0*	0.410 (3.554) [2.823]	-1.059 (1.995) [1.380]	-0.610 (3.118) [2.587]	-2.158 (2.974) [2.511]
500	-1.673 (1.436) [1.103]	-2.075 (1.012) [0.848] 0.747*	-3.435 (1.614) [1.337] 0.787*	4.987 (2.455) [2.039] 0.641*	0.325 (2.198) [1.765] 0.046*	0.407 (1.295) [1.046]	-1.060 (0.598) [0.409]	-0.607 (1.651) [1.290]	-1.921 (1.884) [1.471]
1000	-1.545 (0.847) [0.669]	-2.053 (0.783) [0.619] 0.994*	-3.399 (1.157) [0.899] 0.970*	5.024 (2.069) [1.586] 0.895*	0.309 (1.394) [1.125] 0.083*	0.405 (0.940) [0.737]	-1.045 (0.344) [0.253]	-0.604 (1.006) [0.775]	-1.980 (1.044) [0.843]
1500	-1.595 (0.708) [0.543]	-2.017 (0.630) [0.492] 1*	-3.410 (0.909) [0.679] 0.993*	5.024 (0.895) [0.738] 0.937*	0.306 (1.234) [0.991] 0.089*	0.404 (0.749) [0.606]	-1.035 (0.259) [0.196]	-0.605 (0.866) [0.661]	-1.997 (0.860) [0.679]

**Table 2.** Simulation results for model  $\mathcal{M}_2$ :  $\beta = (1.5, -2.3, 2.5, -3.5, .5)$  (percentage of infected among the susceptibles: 70%).

n	$\hat{\beta}_n$					$\hat{\theta}_n$			
	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$	$\hat{\theta}_{1,n}$	$\hat{\theta}_{2,n}$	$\hat{\theta}_{3,n}$	$\hat{\theta}_{4,n}$
<b>Percentage of immunes = 25%, <math>\theta = (.71, 1, 2, -3)</math></b>									
100	1.512 (1.413) [1.190]	-2.369 (1.165) [0.917] 0.800*	2.522 (1.189) [0.949] 0.814*	-3.518 (2.180) [1.810] 0.271*	0.562 (1.853) [1.582] 0.236*	0.725 (0.979) [0.789]	1.191 (0.830) [0.636]	2.061 (2.728) [2.309]	-2.896 (1.646) [1.264]
500	1.508 (0.905) [0.724]	-2.313 (0.635) [0.478] 0.993*	2.520 (0.568) [0.435] 0.993*	-3.497 (1.198) [0.937] 0.991*	0.514 (0.633) [0.545] 0.629*	0.714 (0.419) [0.335]	1.076 (0.434) [0.268]	2.045 (1.617) [1.290]	-2.959 (0.652) [0.454]
1000	1.499 (0.569) [0.453]	-2.297 (0.488) [0.335] 0.999*	2.508 (0.398) [0.286] 0.998*	-3.502 (0.908) [0.663] 0.997*	0.512 (0.557) [0.479] 0.732*	0.712 (0.308) [0.241]	1.071 (0.387) [0.204]	2.025 (1.178) [0.941]	-2.985 (0.365) [0.273]
1500	1.499 (0.339) [0.331]	-2.299 (0.372) [0.252] 1*	2.497 (0.322) [0.224] 1*	-3.503 (0.701) [0.508] 1*	0.504 (0.522) [0.447] 0.764*	0.708 (0.257) [0.204]	1.050 (0.337) [0.174]	2.012 (0.983) [0.766]	-2.985 (0.289) [0.225]
<b>Percentage of immunes = 50%, <math>\theta = (-.3, -1, 2.1, 1)</math></b>									
100	1.526 (1.887) [1.577]	-2.328 (1.824) [1.507] 0.411*	2.339 (2.170) [1.806] 0.228*	-3.336 (2.570) [2.174] 0.101*	0.488 (2.181) [1.654] 0.060*	-0.332 (0.902) [0.679]	-1.107 (1.204) [0.809]	2.179 (1.946) [1.391]	1.053 (1.567) [1.158]
500	1.517 (0.956) [0.775]	-2.295 (0.772) [0.580] 0.999*	2.635 (0.687) [0.530] 0.963*	-3.397 (1.352) [1.045] 0.924*	0.537 (0.826) [0.648] 0.339*	-0.284 (0.317) [0.255]	-0.983 (0.472) [0.244]	2.127 (0.466) [0.361]	1.041 (1.043) [0.678]
1000	1.517 (0.650) [0.518]	-2.303 (0.531) [0.390] 1*	2.563 (0.467) [0.364] 0.998*	-3.408 (0.962) [0.701] 0.999*	0.512 (0.573) [0.454] 0.399*	-0.296 (0.207) [0.169]	-1.023 (0.157) [0.123]	2.110 (0.310) [0.243]	1.026 (0.532) [0.366]
1500	1.498 (0.473) [0.384]	-2.299 (0.389) [0.281] 1*	2.531 (0.355) [0.281] 1*	-3.365 (0.768) [0.553] 1*	0.513 (0.455) [0.373] 0.451*	-0.295 (0.181) [0.146]	-1.012 (0.122) [0.095]	2.112 (0.260) [0.208]	0.995 (0.339) [0.253]
<b>Percentage of immunes = 75%, <math>\theta = (.4, -1, -.6, -2)</math></b>									
100	1.489 (2.356) [1.934]	-2.384 (2.356) [1.979] 0.109*	2.521 (2.263) [1.944] 0.070*	-3.458 (2.793) [2.389] 0.106*	0.554 (2.326) [1.902] 0.055*	0.420 (1.163) [0.875]	-1.116 (1.602) [1.121]	-0.557 (1.665) [1.232]	-2.123 (1.904) [1.450]
500	1.515 (1.282) [1.047]	-2.381 (1.405) [1.133] 0.929*	2.523 (1.701) [1.317] 0.899*	-3.482 (1.910) [1.527] 0.756*	0.548 (1.757) [1.420] 0.353*	0.418 (0.394) [0.322]	-1.047 (0.396) [0.268]	-0.633 (0.651) [0.482]	-2.084 (0.665) [0.491]
1000	1.490 (0.862) [0.694]	-2.295 (0.916) [0.709] 0.999*	2.522 (1.141) [0.793] 0.999*	-3.515 (1.228) [0.958] 0.929*	0.521 (1.282) [1.013] 0.367*	0.415 (0.243) [0.198]	-1.031 (0.185) [0.142]	-0.627 (0.345) [0.271]	-1.997 (0.398) [0.308]
1500	1.508 (0.721) [0.573]	-2.304 (0.683) [0.505] 1*	2.497 (0.867) [0.575] 1*	-3.512 (0.957) [0.728] 0.973*	0.513 (0.961) [0.765] 0.359*	0.415 (0.221) [0.181]	-1.023 (0.149) [0.119]	-0.613 (0.293) [0.218]	-1.997 (0.299) [0.236]

**Table 3. Simulation results for model  $\mathcal{M}_3$ :**  $\beta = (-1.7, -2.8, 0, -.7, 1.1)$  (percentage of infected among the susceptibles: 30%).

n	$\hat{\beta}_n$					$\hat{\theta}_n$			
	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$	$\hat{\theta}_{1,n}$	$\hat{\theta}_{2,n}$	$\hat{\theta}_{3,n}$	$\hat{\theta}_{4,n}$
<b>Percentage of immunes = 25%, <math>\theta = (.71, 1, 2, -3)</math></b>									
100	-1.732 (1.630) [1.258]	-2.918 (0.631) [0.502] 0.818*	0.017 (1.310) [0.851] 0.033 <sup>†</sup>	-0.764 (1.721) [1.297] 0.053*	1.190 (1.857) [1.569] 0.052*	0.754 (2.681) [2.080]	1.169 (1.686) [1.240]	1.857 (1.983) [1.675]	-2.975 (2.877) [2.325]
500	-1.688 (0.570) [0.441]	-2.908 (0.386) [0.312] 0.998*	-0.016 (0.352) [0.196] 0.055 <sup>†</sup>	-0.751 (0.741) [0.465] 0.301*	1.114 (1.231) [0.956] 0.174*	0.729 (0.837) [0.669]	1.164 (0.634) [0.449]	2.139 (1.494) [1.224]	-3.106 (1.162) [0.898]
1000	-1.698 (0.385) [0.304]	-2.853 (0.272) [0.217] 1*	-0.004 (0.145) [0.115] 0.055 <sup>†</sup>	-0.713 (0.383) [0.301] 0.494*	1.079 (0.781) [0.608] 0.294*	0.726 (0.553) [0.442]	1.080 (0.363) [0.269]	2.083 (1.103) [0.913]	-3.088 (0.771) [0.596]
1500	-1.704 (0.301) [0.234]	-2.837 (0.210) [0.168] 1*	-0.004 (0.126) [0.101] 0.049 <sup>†</sup>	-0.705 (0.277) [0.224] 0.707*	1.110 (0.650) [0.489] 0.479*	0.716 (0.452) [0.359]	1.057 (0.292) [0.205]	2.071 (0.906) [0.743]	-3.087 (0.680) [0.503]
<b>Percentage of immunes = 50%, <math>\theta = (-.3, -1, 2.1, 1)</math></b>									
100	-1.776 (1.879) [1.542]	-2.912 (0.986) [0.807] 0.472*	-0.039 (1.824) [1.404] 0.076 <sup>†</sup>	-0.776 (1.782) [1.493] 0.008*	1.203 (2.056) [1.709] 0.018*	-0.336 (2.904) [2.143]	-1.116 (2.260) [1.533]	1.994 (2.630) [2.053]	1.108 (2.878) [2.129]
500	-1.753 (1.191) [0.919]	-2.918 (0.590) [0.481] 1*	-0.030 (0.490) [0.371] 0.126 <sup>†</sup>	-0.768 (1.361) [1.056] 0.108*	1.194 (1.307) [1.021] 0.196*	-0.279 (0.752) [0.590]	-0.974 (0.806) [0.456]	2.197 (1.312) [0.929]	1.035 (1.053) [0.747]
1000	-1.718 (0.647) [0.525]	-2.853 (0.417) [0.335] 1*	0.005 (0.293) [0.224] 0.084 <sup>†</sup>	-0.719 (0.875) [0.682] 0.148*	1.127 (0.899) [0.710] 0.295*	-0.288 (0.509) [0.414]	-1.003 (0.522) [0.259]	2.149 (0.833) [0.591]	1.021 (0.654) [0.495]
1500	-1.696 (0.551) [0.442]	-2.824 (0.329) [0.259] 0.999*	-0.002 (0.310) [0.181] 0.078 <sup>†</sup>	-0.705 (0.669) [0.517] 0.190*	1.117 (0.701) [0.557] 0.383*	-0.303 (0.387) [0.314]	-1.020 (0.304) [0.186]	2.119 (0.561) [0.423]	1.021 (0.490) [0.385]
<b>Percentage of immunes = 75%, <math>\theta = (.4, -1, -.6, -2)</math></b>									
100	-1.684 (2.086) [1.689]	-2.948 (1.581) [1.313] 0.127*	-0.027 (1.912) [1.591] 0.027 <sup>†</sup>	-0.792 (1.939) [1.621] 0.013*	1.215 (2.648) [2.224] 0.003*	0.497 (3.491) [2.297]	-1.188 (2.037) [1.493]	-0.587 (2.879) [2.183]	-2.120 (2.731) [2.215]
500	-1.774 (1.392) [0.976]	-2.898 (0.908) [0.750] 0.932*	-0.028 (0.993) [0.752] 0.162 <sup>†</sup>	-0.746 (1.651) [1.215] 0.141*	1.197 (1.898) [1.567] 0.084*	0.476 (0.952) [0.720]	-0.923 (1.042) [0.616]	-0.592 (1.396) [0.959]	-1.905 (1.664) [1.156]
1000	-1.745 (0.731) [0.540]	-2.851 (0.631) [0.512] 1*	-0.004 (0.587) [0.435] 0.125 <sup>†</sup>	-0.746 (0.934) [0.715] 0.205*	1.162 (1.430) [1.157] 0.143*	0.473 (0.648) [0.477]	-0.925 (0.797) [0.381]	-0.595 (0.742) [0.545]	-1.927 (0.978) [0.678]
1500	-1.733 (0.596) [0.407]	-2.831 (0.514) [0.410] 1*	-0.004 (0.360) [0.282] 0.114 <sup>†</sup>	-0.746 (0.677) [0.538] 0.247*	1.098 (1.289) [1.002] 0.186*	0.461 (0.513) [0.383]	-0.970 (0.447) [0.247]	-0.595 (0.572) [0.417]	-1.960 (0.845) [0.545]

**Table 4. Simulation results for model  $\mathcal{M}_4$ :**  $\beta = (1.5, -2, 0, 3.5, -4)$  (percentage of infected among the susceptibles: 70%).

n	$\hat{\beta}_n$					$\hat{\theta}_n$			
	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$	$\hat{\theta}_{1,n}$	$\hat{\theta}_{2,n}$	$\hat{\theta}_{3,n}$	$\hat{\theta}_{4,n}$
<b>Percentage of immunes = 25%, <math>\theta = (.71, 1, 2, -3)</math></b>									
100	1.663 (1.209) [1.005]	-2.215 (1.236) [0.975] 0.364*	-0.004 (1.592) [1.125] 0.026†	3.303 (1.775) [1.377] 0.364*	-3.512 (1.895) [1.567] 0.117*	0.793 (1.185) [0.926]	0.939 (1.108) [0.856]	2.216 (1.304) [1.050]	-2.951 (1.989) [1.604]
500	1.580 (0.991) [0.793]	-2.230 (0.959) [0.660] 0.954*	0.039 (0.442) [0.339] 0.137†	3.614 (0.842) [0.679] 0.935*	-3.776 (1.528) [1.108] 0.778*	0.671 (0.486) [0.391]	1.132 (0.745) [0.454]	2.130 (1.041) [0.713]	-2.943 (1.805) [1.324]
1000	1.548 (0.664) [0.519]	-2.096 (0.712) [0.429] 1*	0.028 (0.256) [0.198] 0†	3.559 (0.565) [0.452] 1*	-3.988 (1.098) [0.820] 0.909*	0.685 (0.326) [0.262]	1.051 (0.419) [0.256]	2.081 (0.654) [0.427]	-2.954 (1.241) [0.882]
1500	1.514 (0.560) [0.429]	-2.036 (0.589) [0.314] 0.965*	0.010 (0.206) [0.161] 0.077†	3.538 (0.431) [0.349] 0.965*	-3.997 (0.958) [0.697] 0.958*	0.707 (0.312) [0.246]	1.031 (0.299) [0.192]	2.031 (0.420) [0.318]	-2.975 (0.871) [0.669]
<b>Percentage of immunes = 50%, <math>\theta = (-.3, -1, 2.1, 1)</math></b>									
100	1.472 (1.790) [1.546]	-1.931 (1.663) [1.447] 0.114*	0.022 (2.130) [1.530] 0.184†	3.378 (1.998) [1.691] 0.163*	-3.427 (2.231) [1.946] 0.005*	-0.384 (0.979) [0.690]	-0.913 (1.487) [0.924]	2.286 (1.607) [1.193]	1.276 (1.563) [1.284]
500	1.484 (1.236) [1.020]	-2.034 (1.247) [0.936] 0.878*	0.004 (0.800) [0.518] 0.235†	3.428 (1.180) [0.999] 0.734*	-3.446 (1.826) [1.505] 0.608*	-0.339 (0.303) [0.248]	-0.934 (0.411) [0.232]	2.155 (1.225) [0.633]	1.054 (1.196) [0.945]
1000	1.490 (0.809) [0.652]	-1.956 (0.877) [0.603] 0.994*	0.008 (0.378) [0.282] 0.206†	3.475 (0.960) [0.812] 0.959*	-3.482 (1.565) [1.187] 0.845*	-0.325 (0.206) [0.166]	-0.976 (0.237) [0.130]	2.148 (0.653) [0.401]	1.052 (0.847) [0.658]
1500	1.490 (0.570) [0.458]	-1.987 (0.637) [0.421] 1*	-0.001 (0.322) [0.238] 0.172†	3.492 (0.784) [0.662] 0.989*	-3.763 (0.989) [0.769] 0.892*	-0.308 (0.178) [0.145]	-0.982 (0.208) [0.106]	2.092 (0.475) [0.299]	1.032 (0.649) [0.513]
<b>Percentage of immunes = 75%, <math>\theta = (.4, -1, -.6, -2)</math></b>									
100	1.462 (1.937) [1.643]	-1.936 (2.112) [1.790] 0.042*	-0.012 (2.375) [1.833] 0.143†	3.380 (2.509) [2.211] 0.007*	-3.520 (2.581) [2.146] 0.012*	0.508 (1.050) [0.791]	-1.137 (1.710) [1.018]	-0.710 (1.880) [1.329]	-2.238 (1.705) [1.445]
500	1.456 (1.493) [1.202]	-1.939 (1.418) [1.162] 0.449*	-0.020 (0.820) [0.633] 0.297†	3.453 (1.900) [1.585] 0.116*	-3.466 (2.061) [1.626] 0.231*	0.485 (0.478) [0.372]	-0.933 (0.410) [0.268]	-0.674 (0.644) [0.450]	-2.156 (1.395) [1.156]
1000	1.477 (1.059) [0.851]	-1.947 (1.084) [0.857] 0.912*	0.014 (0.531) [0.419] 0.258†	3.480 (1.521) [1.222] 0.489*	-3.785 (1.601) [1.273] 0.450*	0.462 (0.339) [0.259]	-0.951 (0.300) [0.170]	-0.645 (0.409) [0.291]	-2.051 (1.196) [0.969]
1500	1.482 (0.741) [0.597]	-1.975 (0.731) [0.561] 0.971*	-0.011 (0.368) [0.290] 0.268†	3.492 (1.250) [0.969] 0.692*	-3.801 (1.091) [0.905] 0.558*	0.437 (0.236) [0.190]	-0.961 (0.214) [0.121]	-0.642 (0.345) [0.242]	-2.021 (1.050) [0.838]

**Table 5. Simulation results for the model  $\mathcal{M}_1$ :  $\beta = (-1.7, -2, -3.4, 5, .3)$  if there were no immunes.**

n	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$
100	-1.887 (1.197) [0.938]	-2.474 (1.098) [0.750] 0.985*	-4.126 (1.456) [1.031] 0.996*	5.747 (1.851) [1.437] 0.953*	0.349 (1.825) [1.418] 0.033*
500	-1.749 (0.472) [0.366]	-2.072 (0.282) [0.217] 1*	-3.537 (0.442) [0.332] 1*	5.186 (0.727) [0.558] 1*	0.317 (0.586) [0.469] 0.067*
1000	-1.724 (0.318) [0.253]	-2.027 (0.188) [0.149] 1*	-3.449 (0.275) [0.216] 1*	5.066 (0.456) [0.362] 1*	0.302 (0.436) [0.348] 0.121*
1500	-1.715 (0.253) [0.199]	-2.020 (0.152) [0.121] 1*	-3.437 (0.218) [0.169] 1*	5.053 (0.372) [0.297] 1*	0.298 (0.340) [0.273] 0.145*

**Table 6. Simulation results for the model  $\mathcal{M}_3$ :  $\beta = (-1.7, -2.8, 0, -.7, 1.1)$  if there were no immunes.**

n	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$
100	-1.881 (0.830) [0.659]	-2.937 (0.550) [0.453] 1*	0.002 (0.385) [0.297] 0.044 <sup>†</sup>	-0.753 (0.969) [0.746] 0.137*	1.293 (1.057) [0.803] 0.244*
500	-1.740 (0.348) [0.272]	-2.875 (0.299) [0.228] 1*	-0.002 (0.144) [0.115] 0.047 <sup>†</sup>	-0.718 (0.367) [0.289] 0.540*	1.118 (0.397) [0.314] 0.826*
1000	-1.728 (0.237) [0.188]	-2.823 (0.190) [0.151] 1*	-0.001 (0.095) [0.078] 0.054 <sup>†</sup>	-0.697 (0.243) [0.197] 0.809*	1.116 (0.267) [0.212] 0.989*
1500	-1.711 (0.195) [0.154]	-2.823 (0.159) [0.125] 1*	-0.001 (0.081) [0.065] 0.047 <sup>†</sup>	-0.702 (0.202) [0.162] 0.935*	1.104 (0.229) [0.184] 0.998*



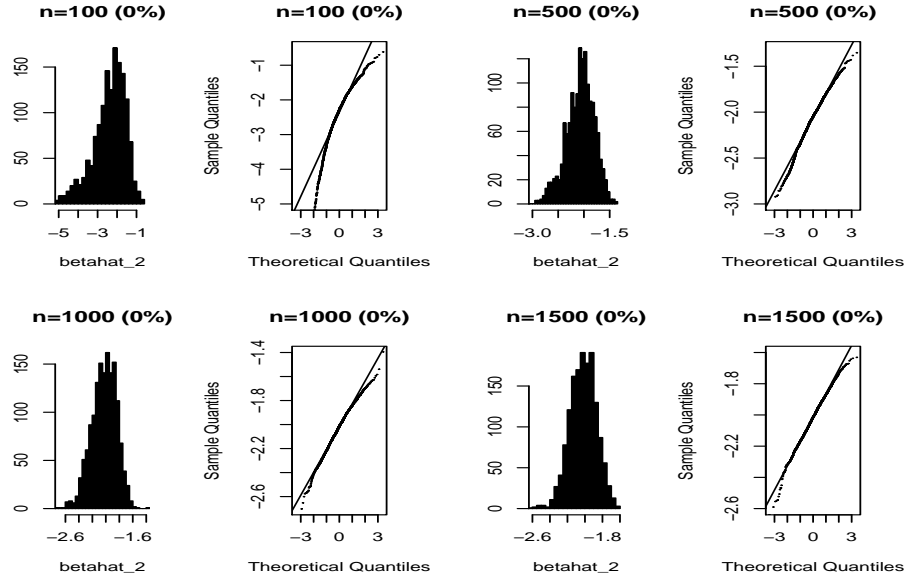
**Table 7. Results of a "naive" analysis of model  $\mathcal{M}_1$ :  $\beta = (-1.7, -2, -3.4, 5, .3)$ .**

n	$\hat{\beta}_n$				
	$\hat{\beta}_{1,n}$	$\hat{\beta}_{2,n}$	$\hat{\beta}_{3,n}$	$\hat{\beta}_{4,n}$	$\hat{\beta}_{5,n}$
<b>Percentage of immunes = 25%</b>					
100	-1.093	-0.376	0.115	2.358	-1.288
	(3.448)	(3.792)	(4.844)	(5.860)	(4.619)
	[2.082]	[2.048]	[3.721]	[5.131]	[2.526]
		0.032*	0.294*	0.171*	0.497*
500	-1.376	-0.085	0.265	1.815	-1.572
	(3.609)	(1.929)	(4.058)	(5.159)	(2.845)
	[2.179]	[1.917]	[3.751]	[4.657]	[2.227]
		0.053*	0.559*	0.446*	0.701*
1000	-1.290	-0.158	0.158	2.410	-0.966
	(2.944)	(2.182)	(3.565)	(4.871)	(2.760)
	[2.019]	[1.921]	[3.558]	[4.304]	[1.843]
		0.046*	0.624*	0.535*	0.751*
1500	-1.171	-0.112	0.162	1.962	-1.044
	(2.138)	(1.902)	(3.570)	(4.885)	(2.659)
	[1.867]	[1.890]	[3.562]	[4.437]	[1.847]
		0.052*	0.652*	0.570*	0.778*
<b>Percentage of immunes = 50%</b>					
100	-1.953	-0.228	-1.048	3.165	0.171
	(3.843)	(4.118)	(5.234)	(5.194)	(0.506)
	[1.745]	[2.181]	[3.101]	[4.347]	[0.397]
		0.028*	0.508*	0.437*	0.139*
500	-1.611	-0.237	-0.663	2.754	0.426
	(1.583)	(2.603)	(3.193)	(5.193)	(0.458)
	[1.231]	[1.902]	[2.845]	[4.383]	[0.392]
		0.046*	0.704*	0.658*	0.457*
1000	-1.755	-0.148	-0.617	1.274	0.432
	(1.348)	(1.871)	(2.836)	(4.088)	(0.434)
	[1.175]	[1.855]	[2.800]	[3.725]	[0.371]
		0.043*	0.759*	0.731*	0.587*
1500	-1.446	-0.124	-0.621	1.425	0.487
	(1.281)	(1.891)	(2.800)	(3.632)	(0.433)
	[1.143]	[1.878]	[2.778]	[3.574]	[0.378]
		0.051*	0.799*	0.755*	0.632*
<b>Percentage of immunes = 75%</b>					
100	-1.665	-0.284	-1.095	1.434	0.178
	(4.897)	(4.546)	(5.127)	(6.948)	(6.963)
	[2.514]	[2.293]	[2.981]	[5.974]	[2.119]
		0.037*	0.515*	0.083*	0.237*
500	-1.746	-0.484	-1.032	0.751	0.446
	(3.847)	(3.748)	(4.757)	(6.987)	(6.769)
	[2.091]	[2.028]	[2.936]	[5.969]	[1.901]
		0.041*	0.696*	0.305*	0.596*
1000	-1.520	-0.261	-0.745	0.252	0.321
	(2.883)	(4.496)	(2.679)	(5.839)	(3.468)
	[1.857]	[2.110]	[2.655]	[5.662]	[1.194]
		0.038*	0.788*	0.481*	0.696*
1500	-1.510	-0.120	-0.757	0.132	0.315
	(2.550)	(1.902)	(2.671)	(6.016)	(3.658)
	[1.791]	[1.887]	[2.649]	[5.659]	[1.427]
		0.041*	0.801*	0.567*	0.727*

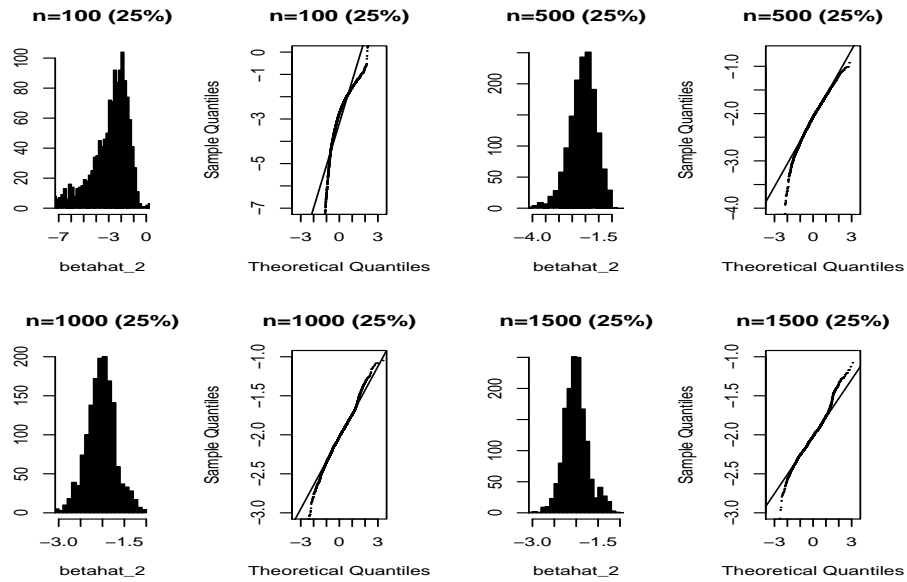
**Table 8.** Estimated probability  $p(\mathbf{x})$  for the models  $\mathcal{M}_1$  ( $p(\mathbf{x}) = 0.250$ ) and  $\mathcal{M}_2$  ( $p(\mathbf{x}) = 0.343$ ).

	0% of immune		25% of immune		50% of immune		75% of immune	
n	$\mathcal{M}_1$	$\mathcal{M}_2$	$\mathcal{M}_1$	$\mathcal{M}_2$	$\mathcal{M}_1$	$\mathcal{M}_2$	$\mathcal{M}_1$	$\mathcal{M}_2$
100	0.229 (0.118) [0.097] 0.964* 0.407 $\mp$	0.335 (0.117) [0.095] 0.959* 0.424 $\mp$	0.267 (0.191) [0.149] 0.964* 0.596 $\mp$	0.363 (0.164) [0.132] 0.978* 0.523 $\mp$	0.264 (0.264) [0.202] 0.987* 0.719 $\mp$	0.353 (0.215) [0.174] 0.949* 0.627 $\mp$	0.266 (0.360) [0.293] 0.951* 0.859 $\mp$	0.362 (0.367) [0.318] 0.831* 0.756 $\mp$
500	0.246 (0.046) [0.037] 0.957* 0.182 $\mp$	0.343 (0.049) [0.039] 0.959* 0.192 $\mp$	0.261 (0.079) [0.056] 0.930* 0.216 $\mp$	0.356 (0.069) [0.054] 0.898* 0.228 $\mp$	0.263 (0.110) [0.083] 0.921* 0.343 $\mp$	0.358 (0.084) [0.068] 0.943* 0.319 $\mp$	0.255 (0.188) [0.150] 0.891* 0.572 $\mp$	0.354 (0.201) [0.160] 0.701* 0.482 $\mp$
1000	0.247 (0.034) [0.027] 0.948* 0.128 $\mp$	0.342 (0.035) [0.028] 0.954* 0.136 $\mp$	0.255 (0.051) [0.037] 0.887* 0.149 $\mp$	0.352 (0.052) [0.039] 0.875* 0.164 $\mp$	0.258 (0.071) [0.055] 0.931* 0.248 $\mp$	0.351 (0.061) [0.049] 0.932* 0.226 $\mp$	0.254 (0.137) [0.106] 0.894* 0.436 $\mp$	0.350 (0.134) [0.106] 0.651* 0.359 $\mp$
1500	0.249 (0.028) [0.022] 0.945* 0.105 $\mp$	0.343 (0.028) [0.023] 0.958* 0.112 $\mp$	0.252 (0.037) [0.028] 0.907* 0.119 $\mp$	0.348 (0.038) [0.028] 0.877* 0.133 $\mp$	0.253 (0.062) [0.044] 0.932* 0.199 $\mp$	0.352 (0.047) [0.038] 0.929* 0.185 $\mp$	0.251 (0.106) [0.085] 0.904* 0.373 $\mp$	0.348 (0.108) [0.084] 0.641* 0.301 $\mp$

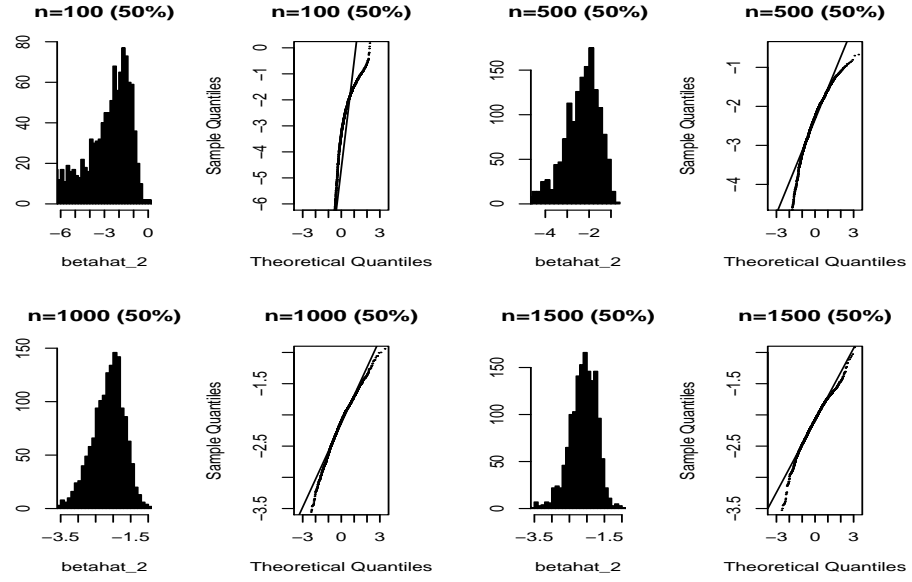
Note: ( $\cdot$ ): root mean square error. [ $\cdot$ ]: mean absolute error. \*: empirical coverage probability.  $\mp$ : average length of confidence intervals. For each percentage of immunes, the percentages of infected among the susceptibles are respectively 30% ( $\mathcal{M}_1$ ) and 70% ( $\mathcal{M}_2$ ). All results are based on 1500 replicates.



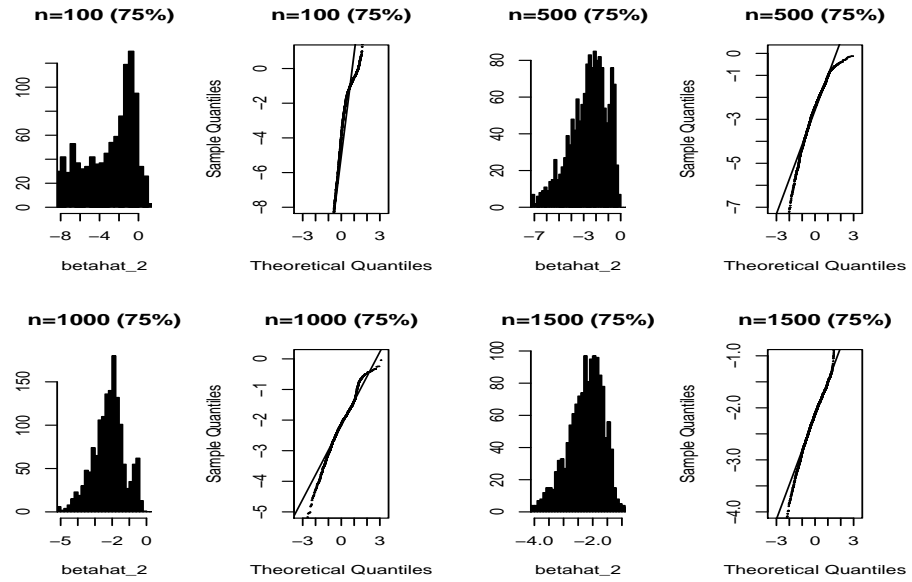
**Figure 1.** Histograms and Q-Q plots for  $\hat{\beta}_{2,n}$  in model  $\mathcal{M}_1$ , with no immunes in the sample (the percentage of immunes is given in brackets).  $n$  is the sample size. All results are based on 1500 simulated datasets.



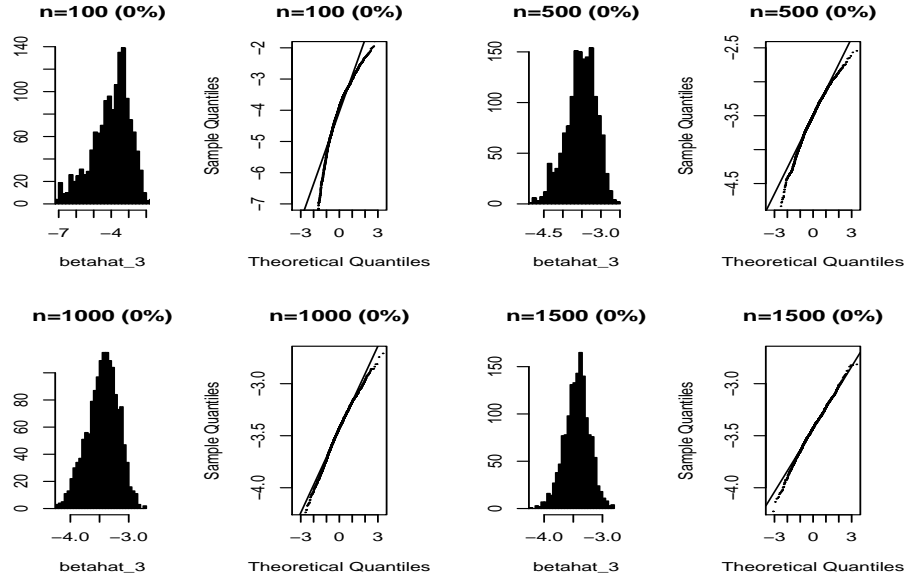
**Figure 2.** Histograms and Q-Q plots for  $\hat{\beta}_{2,n}$  in model  $\mathcal{M}_1$ , with 25% of immunes.



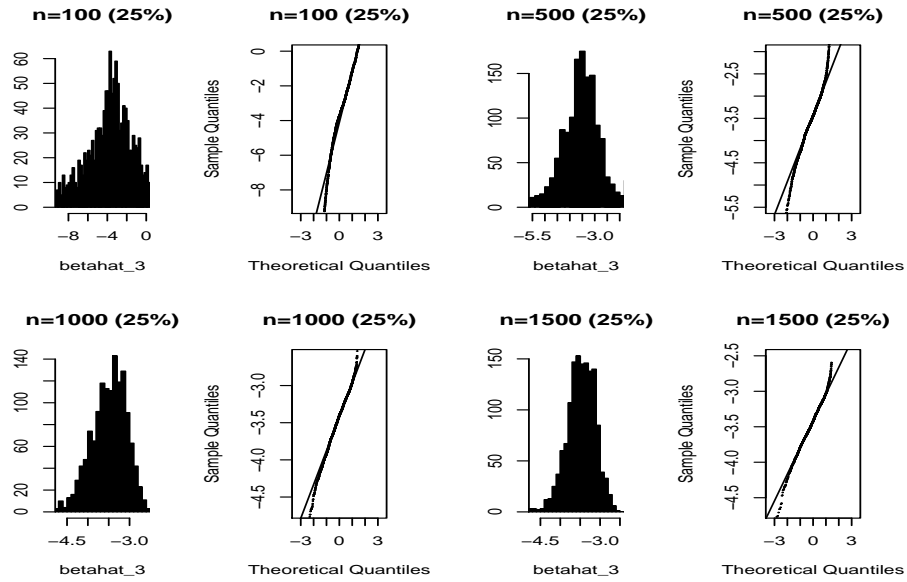
**Figure 3.** Histograms and Q-Q plots for  $\hat{\beta}_{2,n}$  in model  $\mathcal{M}_1$ , with 50% of immunes.



**Figure 4.** Histograms and Q-Q plots for  $\hat{\beta}_{2,n}$  in model  $\mathcal{M}_1$ , with 75% of immunes.



**Figure 5.** Histograms and Q-Q plots for  $\hat{\beta}_{3,n}$  in model  $\mathcal{M}_1$ , with no immunes in the sample.



**Figure 6.** Histograms and Q-Q plots for  $\hat{\beta}_{3,n}$  in model  $\mathcal{M}_1$ , with 25% of immunes.

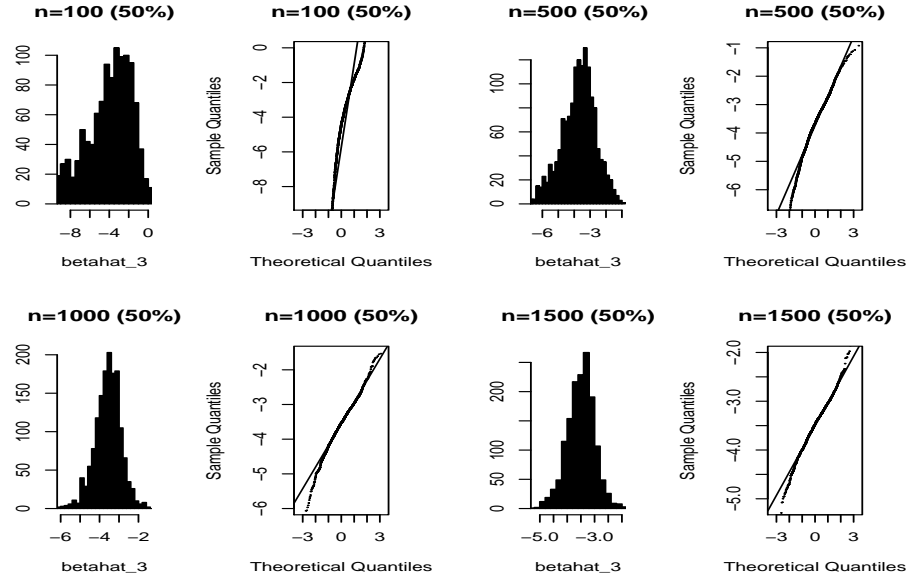


Figure 7. Histograms and Q-Q plots for  $\hat{\beta}_{3,n}$  in model  $\mathcal{M}_1$ , with 50% of immunes.

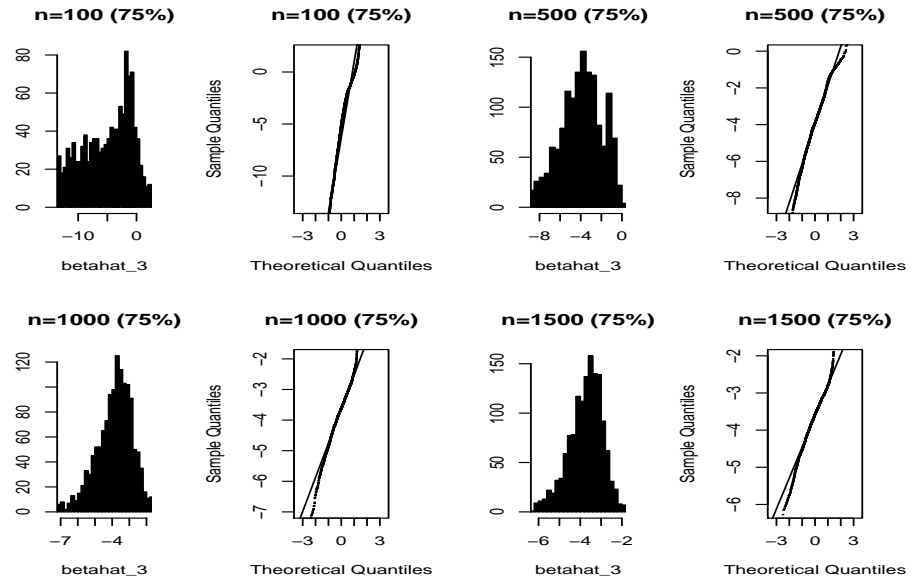
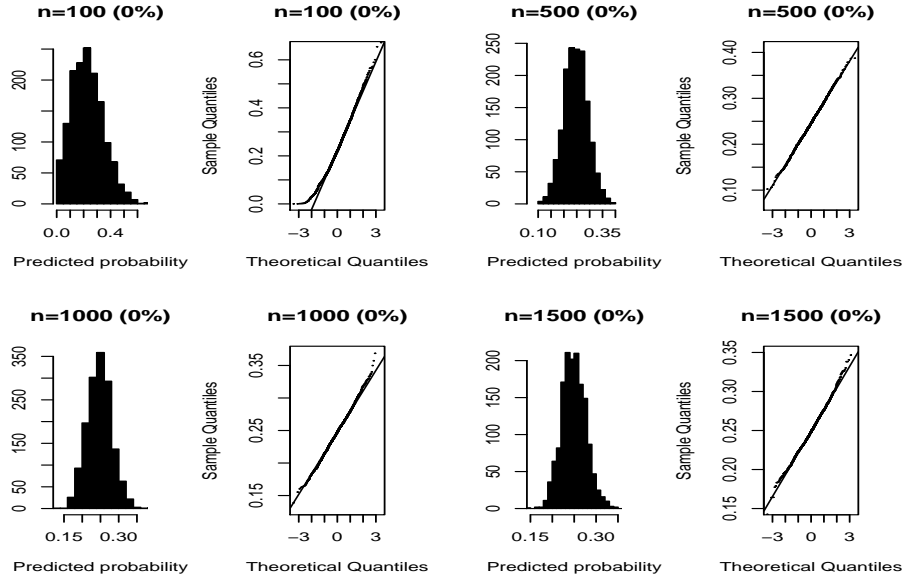
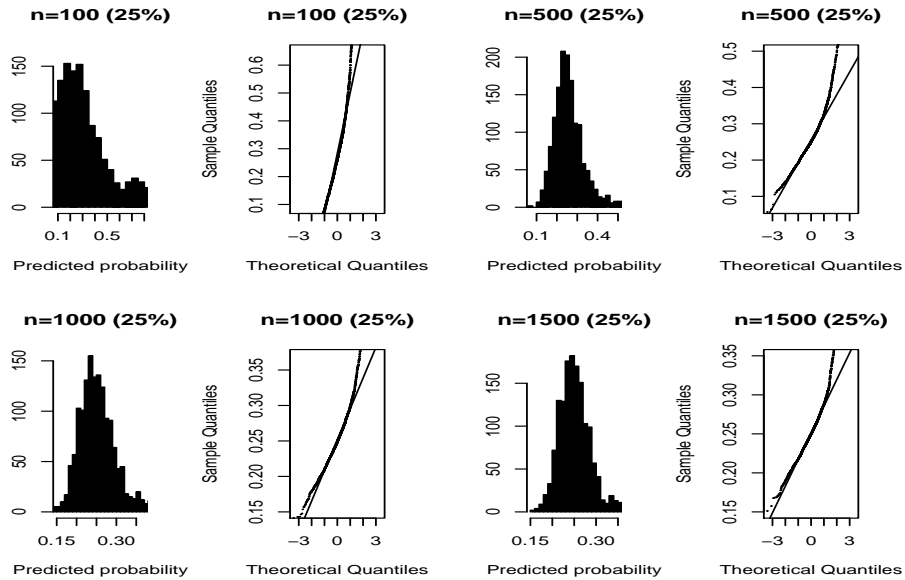


Figure 8. Histograms and Q-Q plots for  $\hat{\beta}_{3,n}$  in model  $\mathcal{M}_1$ , with 75% of immunes.



**Figure 9.** Histograms and Q-Q plots for  $\hat{p}_n(x)$  in model  $\mathcal{M}_1$ , with no immunes in the sample.



**Figure 10.** Histograms and Q-Q plots for  $\hat{p}_n(x)$  in model  $\mathcal{M}_1$ , with 25% of immunes.

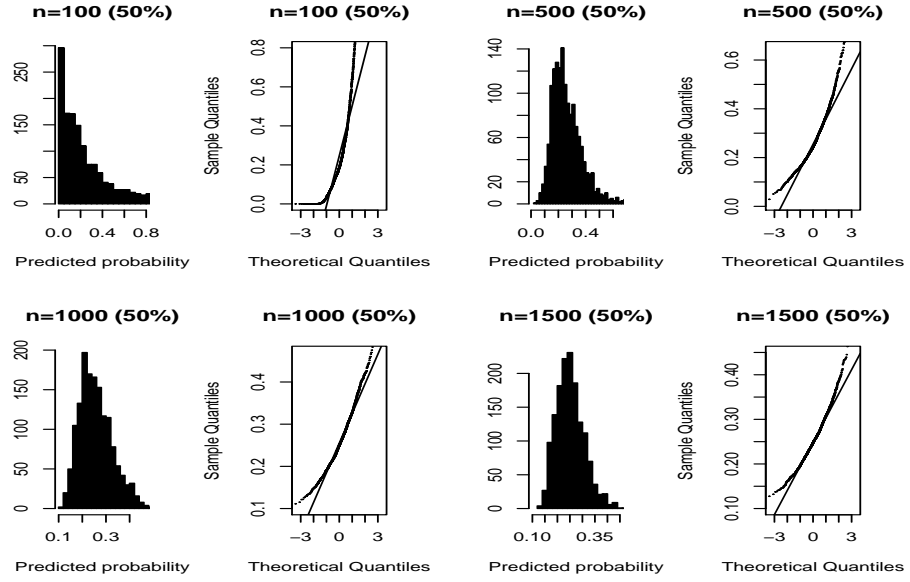


Figure 11. Histograms and Q-Q plots for  $\hat{p}_n(x)$  in model  $\mathcal{M}_1$ , with 50% of immunes.

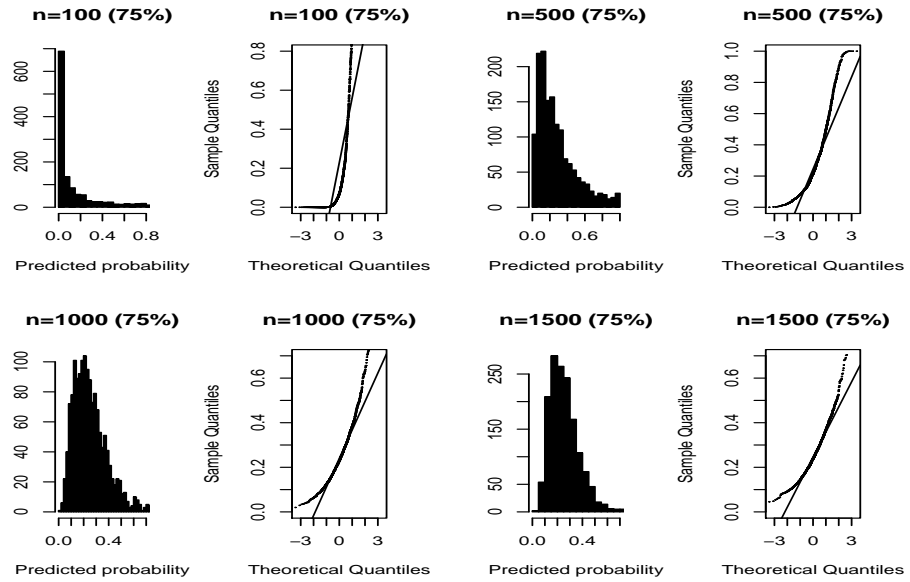


Figure 12. Histograms and Q-Q plots for  $\hat{p}_n(x)$  in model  $\mathcal{M}_1$ , with 75% of immunes.